

# On the Geometry of Bayesian Inference

Miguel DE CARVALHO, Garritt L. PAGE, and Bradley J. BARNEY

## Abstract

We provide a geometric interpretation to Bayesian inference that allows us to introduce a natural measure of the level of agreement between priors, likelihoods, and posteriors. The starting point for the construction of our geometry is the simple observation that the marginal likelihood can be regarded as an inner product between the prior and the likelihood. A key concept in our geometry is that of compatibility, a measure which is based on the same construction principles as Pearson correlation, but which can be used to assess how much the prior agrees with the likelihood, to gauge the sensitivity of the posterior to the prior, and to quantify the coherency of the opinions of two experts. Estimators for all the quantities involved in our geometric setup are discussed, which can be directly computed from the posterior simulation output. Some examples are used to illustrate our methods, including data related to on-the-job drug usage, midge wing length, and prostate cancer.

KEYWORDS: Bayesian inference; Geometry; Harmonic mean estimator; Hilbert spaces; Marginal likelihood; Normalizing constant; Prior-data conflict.

## 1 Introduction

The increased complexity of models posed in fields such as biology, ecology, and epidemiology (to name a few) has led many practitioners to adopt Bayesian methodologies. This trend is not necessarily motivated by philosophical underpinnings, rather no alternative machinery capable of fitting posed models exists. Thus, in a way, Bayesian methods have become more or less mainstream, and this has led to an increased need for model assessment metrics that are quickly calculated and easily interpreted. A welcome metric to practitioners would be one that is able to guide decisions in the

---

Miguel de Carvalho is Lecturer in Statistics, School of Mathematics, The University of Edinburgh, UK (*e-mail: miguel.decarvalho@ed.ac.uk*). Garritt L. Page is Assistant Professor of Statistics, Department of Statistics, Brigham Young University, Provo, Utah (*e-mail: page@stat.byu.edu*). Bradley J. Barney is Visiting Assistant Professor of Statistics, Department of Statistics, Brigham Young University, Provo, Utah (*e-mail: barney@stat.byu.edu*). The authors thank José Quinlan for research assistantship and for numerous discussions, and extend their thanks to Vanda Inácio de Carvalho, Anthony Davison, Duvan Henao, Wesley Johnson, Antónia Turkman, and Feridun Turkman for constructive comments and discussions. The research was partially supported by the Fondecyt projects no. 11121186 and 11121131.

model building process by providing a quick assessment of the level of agreement or influence that each component of Bayes theorem has on inference and predictions.

Assessing the influence that prior distributions and/or likelihoods have on posterior inference has been a topic of research for some time. One commonly used ad-hoc method suggests simply fitting a Bayes model using a few competing priors, then visually (or numerically) assessing changes in the posterior as a whole or using some pre-specified posterior summary. More rigorous approaches have also been developed. [Lavine \(1991\)](#) developed a framework to assess sensitivity of posterior inference to sampling distribution (likelihood) and the priors. [Berger \(1991\)](#) introduced the concept of Bayesian robustness which includes perturbation models (see also [Berger and Berliner 1986](#)). More recently, [Evans and Jang \(2011\)](#) have compared information available in two competing priors. Related to this work, [Gelman et al. \(2011\)](#) advocates the use of so-called weakly informative priors that purposely incorporate less information than available as a means of regularizing. Work has also been dedicated to the so-called prior-data conflict which aims to assess the level of agreement between prior and likelihood (see [Evans and Moshonov 2006](#), [Walter and Augustin 2009](#), [Al Labadi and Evans 2016](#)). Such conflict can be of interest in a wealth of situations, such as for assessing how much an expert agrees with the data, or for evaluating how much prior and likelihood information are at odds at the node level in a hierarchical model (see [Scheel, Green and Rougier, 2011](#), and references therein). Regarding sensitivity of the posterior distribution to prior specifications, [Lopes and Tobias \(2011\)](#) provide a fairly accessible overview.

We argue that a geometric representation of the prior, likelihood, and posterior distribution encourages understanding of their interplay. Considering Bayes methodologies from a geometric perspective is not new, but none of the existing geometric perspectives has been designed with the goal of providing a summary on the agreement or impact that each component of Bayes theorem has on inference and predictions. [Aitchison \(1971\)](#) used a geometric perspective to build intuition behind each component of Bayes theorem. [Zhu, Ibrahim and Tang \(2011\)](#) defined a manifold on which a Bayesian perturbation analysis can be carried out by perturbing data, prior and likelihood simultaneously, [Shortle and Mendel \(1996\)](#) used a geometric approach to draw conditional distributions in arbitrary coordinate systems, and [Agarawal and Daumé \(2010\)](#) argued that conjugate priors of posterior distributions belong to the same geometry giving an appealing interpretation of hyperparameters.

The novel contribution we aim to make here is the development of easily computed metrics that provide an informative preliminary ‘snap-shot’ regarding comparisons between prior and likelihood (to assess the level of agreement between prior and data), prior and posterior (to determine the influence that prior has on inference), and prior versus prior (to compare ‘informativeness’—i.e., a density’s

peakedness—and/or congruence of two competing priors). To this end, we treat each component of Bayes theorem as an element of a geometry formally constructed using concepts from Hilbert spaces and tools from abstract geometry. Because of this, it is possible to calculate norms, inner products, and angles between vectors. Not only do each of these numeric summaries have intuitively appealing individual interpretations, but they may also be combined to construct a unitless measure of compatibility, which can be used to assess how much the prior agrees with the likelihood, to gauge the sensitivity of the posterior to the prior, and to quantify the coherency of the opinions of two experts. Further, estimating our measures of ‘similarity’ is straightforward and can actually be carried out within an MCMC algorithm as is typically employed in a Bayesian analysis.

To facilitate the illustration of ideas, concepts, and methods we reference the following simple example (which is found in [Christensen et al. 2011](#), pp. 26–27) through the first few sections of this article.

#### ON-THE-JOB DRUG USAGE TOY EXAMPLE

Suppose interest lies in estimating the proportion  $\theta \in [0, 1]$  of US transportation industry workers that use drugs on the job. Suppose  $n = 10$  workers were selected and tested with the 2nd and 7th testing positive. Let  $\mathbf{y} = (Y_1, \dots, Y_{10})$  with  $Y_i = 1$  denoting that the  $i$ th worker tested positive and  $Y_i = 0$  otherwise. A natural data model for these data would be  $\mathbf{y} \mid \theta \stackrel{\text{iid}}{\sim} \text{Bern}(\theta)$ . A prior distribution that is typically used in this situation is  $\theta \sim \text{Beta}(a, b)$  for  $a, b > 0$  and produces the following posterior distribution  $\theta \mid \mathbf{y} \sim \text{Beta}(a^*, b^*)$  with  $a^* = n_1 + a$  and  $b^* = n - n_1 + b$ .

Some natural questions one may ask and that we aim to quantify are: How compatible is the likelihood with this prior choice? How similar are the posterior and prior distributions? How does the choice of  $\text{Beta}(a, b)$  compare to other possible prior distributions? We provide a unified treatment to answer the questions above. While the drug usage example provides a recurring backdrop that we consistently call upon, additional examples are used throughout the paper to illustrate our methods.

The rest of the article is organized as follows. Section 2 introduces the basic geometric framework in which we work and provides definitions and interpretations of norms and inner-products. Section 3 generalizes how Bayes theorem employs a likelihood to recast a prior density to obtain a posterior density. Section 4 contains computational details. Section 5 provides a regression example illustrating utility of our metric. Section 6 conveys some concluding remarks. Proofs are given in the Appendix.

## 2 Bayes geometry

### 2.1 A geometric view of Bayes theorem

Suppose the inference of interest is over a parameter  $\boldsymbol{\theta}$  which takes values on  $\Theta \subseteq \mathbb{R}^p$ . We consider the space of square integrable functions  $L_2(\Theta)$ , and use the geometry of the Hilbert space  $\mathcal{H} = (L_2(\Theta), \langle \cdot, \cdot \rangle)$ , with inner-product

$$\langle g, h \rangle = \int_{\Theta} g(\boldsymbol{\theta})h(\boldsymbol{\theta}) \, \mathrm{d}\boldsymbol{\theta}, \quad g, h \in L_2(\Theta). \quad (1)$$

The fact that  $\mathcal{H}$  is an Hilbert space is often known in mathematical parlance as the Riesz–Fischer theorem; for a proof see [Cheney \(2001, p. 411\)](#). Borrowing geometric terminology from linear spaces, we refer to the elements of  $L_2(\Theta)$  as vectors, and assess their ‘magnitudes’ through the use of the norm induced by the inner product in (1), i.e.,  $\|\cdot\| = (\langle \cdot, \cdot \rangle)^{1/2}$ .

The starting point for constructing our geometry is the observation that Bayes theorem can be written using the inner-product in (1) as follows

$$p(\boldsymbol{\theta} \mid \mathbf{y}) = \frac{\pi(\boldsymbol{\theta})f(\mathbf{y} \mid \boldsymbol{\theta})}{\int_{\Theta} \pi(\boldsymbol{\theta})f(\mathbf{y} \mid \boldsymbol{\theta}) \, \mathrm{d}\boldsymbol{\theta}} = \frac{\pi(\boldsymbol{\theta})\ell(\boldsymbol{\theta})}{\langle \pi, \ell \rangle}, \quad (2)$$

where  $\ell(\boldsymbol{\theta}) = f(\mathbf{y} \mid \boldsymbol{\theta})$  denotes the likelihood,  $\pi(\boldsymbol{\theta})$  is a prior density,  $p(\boldsymbol{\theta} \mid \mathbf{y})$  is the posterior density and  $\langle \pi, \ell \rangle = \int_{\Theta} f(\mathbf{y} \mid \boldsymbol{\theta})\pi(\boldsymbol{\theta}) \, \mathrm{d}\boldsymbol{\theta}$  is the so-called marginal likelihood or integrated likelihood. The inner product in (1) naturally leads to considering  $\pi$  and  $\ell$  that are in  $L_2(\Theta)$ , which is compatible with a wealth of parametric models and proper priors. By considering  $p$ ,  $\pi$ , and  $\ell$  as vectors with different magnitudes and directions, Bayes theorem simply indicates how one might recast the prior vector so to obtain the posterior vector. The likelihood vector is used to enlarge/reduce the magnitude and suitably tilt the direction of the prior vector in a sense that will be made precise below.

The marginal likelihood  $\langle \pi, \ell \rangle$  is simply the inner product between the likelihood and the prior, and hence can be understood as a natural measure of agreement between the prior and the likelihood. To make this more concrete, define the *angle measure* between the prior and the likelihood as

$$\pi \angle \ell = \arccos \frac{\langle \pi, \ell \rangle}{\|\pi\| \|\ell\|}. \quad (3)$$

Since  $\pi$  and  $\ell$  are nonnegative, the angle between the prior and the likelihood can only be acute or right, i.e.,  $\pi \angle \ell \in [0, 90^\circ]$ . The closer  $\pi \angle \ell$  is to  $0^\circ$ , the greater the agreement between the prior and the likelihood. Conversely, the closer  $\pi \angle \ell$  is to  $90^\circ$ , the greater the disagreement between prior and likelihood. In the pathological case where  $\pi \angle \ell = 90^\circ$  (which requires the prior and the likelihood to have all of their mass on disjoint sets), we say that the prior is orthogonal to the likelihood. Bayes

theorem is incompatible with a prior being completely orthogonal to the likelihood as  $\pi \angle \ell = 90^\circ$  indicates that  $\langle \pi, \ell \rangle = 0$ , thus leading to a division by zero in (2). Similar to the correlation coefficient for random variables in  $L_2(\Omega, \mathbb{B}_\Omega, P)$ —with  $\mathbb{B}_\Omega$  denoting the Borel sigma-algebra over the sample space  $\Omega$ —, our target object of interest is given by a standardized inner product

$$\kappa_{\pi, \ell} = \frac{\langle \pi, \ell \rangle}{\|\pi\| \|\ell\|}. \quad (4)$$

The quantity  $\kappa_{\pi, \ell}$  quantifies how much an expert’s opinion agrees with the data, thus providing a natural measure of prior-data compatibility. As it will become apparent,  $\kappa_{\pi, \ell}$  can be used as a simple alternative to the procedure developed by [Evans and Moshonov \(2006\)](#) to assess prior-data conflict (see also [Walter and Augustin 2009](#); [Scheel, Green and Rougier 2011](#); [Al Labadi and Evans 2016](#)), but with some differences that will be highlighted in the discussion below about Figure 1.

Although not with the express interest of making comparisons between two functions, it should be noted that the idea of angles between functions has appeared in the context of functional data analysis ([Ramsey and Silverman, 1997](#), p. 6).

Before exploring (4) more fully by providing interpretations and properties we concretely define how the term ‘geometry’ will be used throughout the paper. The following definition of abstract geometry can be found in [Millman and Parker \(1991](#), p. 17).

**Definition 1** (Abstract geometry). *An abstract geometry  $\mathcal{A}$  consists of a pair  $\{\mathcal{P}, \mathcal{L}\}$ , where the elements of set  $\mathcal{P}$  are designed as points, and the elements of the collection  $\mathcal{L}$  are designed as lines, such that:*

1. *For every two points  $A, B \in \mathcal{P}$ , there is a line  $l \in \mathcal{L}$ .*
2. *Every line has at least two points.*

Our abstract geometry of interest is  $\mathcal{A} = \{\mathcal{P}, \mathcal{L}\}$ , where  $\mathcal{P} = L_2(\Theta)$  and the set of all lines is

$$\mathcal{L} = \{g + kh : g, h \in L_2(\Theta)\}.$$

Hence, in our setting points can be, for example, prior densities, posterior densities, or likelihoods, as long as they are in  $L_2(\Theta)$ . Lines are elements of  $\mathcal{L}$ , as defined in (2.1), so that for example if  $g$  and  $h$  are densities, line segments in our geometry consist of all possible mixture distributions which can be obtained from  $g$  and  $h$ , i.e.,

$$\{\lambda g + (1 - \lambda)h : \lambda \in [0, 1]\}.$$

Vectors in  $\mathcal{A} = \{\mathcal{P}, \mathcal{L}\}$  are defined through the difference of elements in  $\mathcal{P} = L_2(\Theta)$ . For example, let  $g \in L_2(\Theta)$  and let  $0 \in L_2(\Theta)$ . Then  $g = g - 0 \in L_2(\Theta)$ , and hence  $g$  can be regarded both as a point

and as a vector. If  $g, h \in L_2(\Theta)$  are vectors then we say that  $g$  and  $h$  are collinear if there exists  $k \in \mathbb{R}$ , such that  $g(\theta) = kh(\theta)$ . Put differently, we say  $g$  and  $h$  are collinear if  $g(\theta) \propto h(\theta)$ , for all  $\theta \in \Theta$ .

For any two points in the geometry under consideration, we define their compatibility as a standardized inner product (with (4) being a particular case).

**Definition 2** (Compatibility). *The compatibility between points in the geometry under consideration is the mapping  $\kappa : L_2(\Theta) \times L_2(\Theta) \rightarrow [0, 1]$  defined as*

$$\kappa_{g,h} = \frac{\langle g, h \rangle}{\|g\| \|h\|}, \quad g, h \in L_2(\Theta). \quad (5)$$

The concept of compatibility in Definition 2 is based on the same construction principles as the Pearson correlation coefficient, which would be based however on the inner product

$$\langle X, Y \rangle = \int_{\Omega} XY \, dP, \quad X, Y \in L_2(\Omega, \mathbb{B}_{\Omega}, P), \quad (6)$$

instead of the inner product in (1). For a few selected  $\pi$ 's,  $\kappa_{\pi,p}$  can be used to gauge the sensitivity of the posterior to the prior specification. Also,  $\kappa_{\pi_1, \pi_2}$  might quantify the compatibility of different priors, and hence it can be used to assess the coherency of the opinions of two experts. As an illustration consider the following simple example.

**Example 1.** Consider the following densities  $\pi_0(\theta) = I_{(0,1)}(\theta)$ ,  $\pi_1(\theta) = I_{(0,2)}(\theta)$ ,  $\pi_2(\theta) = I_{(1,2)}(\theta)$ , and  $\pi_3(\theta) = I_{(1,3)}(\theta)$ . Note that  $\|\pi_0\| = \|\pi_2\| = 1$ ,  $\|\pi_1\| = \|\pi_3\| = \sqrt{2}/2$ , and; further,  $\kappa_{\pi_0, \pi_1} = \kappa_{\pi_2, \pi_3} = \sqrt{2}/2$ , thus implying that  $\pi_0 \angle \pi_1 = \pi_2 \angle \pi_3 = 45^\circ$ . Note further that  $\kappa_{\pi_0, \pi_2} = 0$ , and hence  $\pi_0 \perp \pi_2$ .

As can be observed in Example 1,  $(\pi_a \angle \pi_b)/90^\circ$  is a natural measure of distinctiveness of two densities. In addition, Example 1 shows us how different distributions can be associated to the same norm and angle. Hence, as expected, any Cartesian representation  $(x, y) \mapsto (\|\cdot\| \cos(\cdot \angle \cdot), \|\cdot\| \sin(\cdot \angle \cdot))$ , will only allow us to represent some features of the corresponding distributions, but will not allow us to identify the distributions themselves.

To begin building intuition regarding the values produced by  $\kappa_{\pi, \ell}$ , we provide Figure 1. In the figure,  $\ell$  is set to  $N(0, 1)$  while  $\pi = N(m, \sigma)$  varies according to  $m$  and  $\sigma$ . The left plot corresponds to fixing  $\sigma = 1$  and varying  $m$  while in the right plot  $m = 0$  is fixed and  $\sigma$  varies. Notice that in plot (i)  $\kappa_{\pi, \ell} = 0.1$  corresponds to distributions whose means are approximately 3 standard deviations apart while a  $\kappa_{\pi, \ell} = 0.9$  corresponds to distributions whose means are approximately 0.65 standard deviations apart. Connecting specific values of  $\kappa$  to specific standard deviation distances between means seems like a natural way to quickly get a rough idea of relative differences between two distributions. In plot (ii) it appears that if both distributions are centered at the same value, then one distribution

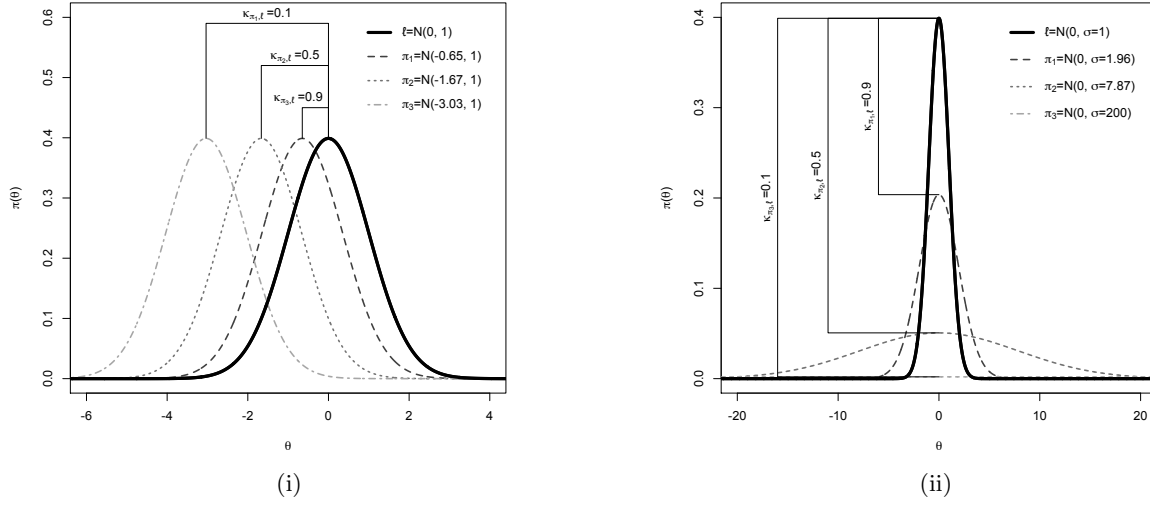


Figure 1: Values of  $\kappa_{\pi, \ell}$  when both  $\pi$  and  $\ell$  are both Gaussian distributions. Figure (i) depicts Gaussian distributions whose means become more separated, while Figure (ii) depicts Gaussian distributions that become progressively more diffuse.

must be very disperse relative to the other to produce  $\kappa$  values that are small (e.g.,  $\leq 0.1$ ). This makes sense as there always exists some mass intersection between the two distributions considered. In this scenario—an especially diffuse prior—[Evans and Moshonov 2006](#) would conclude that no prior-data conflict exists. Their method relies on determining how extreme the data are relative to the induced prior-predictive distribution, and an increasingly disperse prior would make any given data set seem increasingly less extreme. We consider  $\kappa_{\pi, \ell}$  more a measure of prior-data *compatibility* than prior-data *conflict* inasmuch as it can be small not only when there are differences in the locations of the prior and likelihood but also when there are differences in the peakedness of the distributions. Some further comments regarding our geometry are in order:

- Two different densities  $\pi_1$  and  $\pi_2$  cannot be collinear: If  $\pi_1 = k\pi_2$ , then  $k = 1$ , otherwise  $\int \pi_2(\theta) d\theta \neq 1$ .
- A density can be collinear to a likelihood: If the prior is Uniform  $p(\theta | \mathbf{y}) \propto \ell(\theta)$ , and hence the posterior is collinear to the likelihood, i.e., in such case the posterior simply consists of a renormalization of the likelihood.
- Our geometry is compatible with having two likelihoods be collinear, and thus it can be used to rethink the strong likelihood principle ([Berger and Wolpert, 1988](#)). Let  $\ell$  and  $\ell^*$  be the likelihoods based on observing  $\mathbf{y}$  and  $\mathbf{y}^*$ , respectively. The strong likelihood principle states that if  $\ell(\theta) = f(\theta | \mathbf{y}) \propto f(\theta | \mathbf{y}^*) = \ell^*(\theta)$ , then the *same* inference should be drawn from both

samples. According to our geometry, this would mean that likelihoods with the same direction should yield the same inference.

## 2.2 Norms and their interpretation

As  $\kappa_{\pi,\ell}$  is comprised of function norms, we dedicate some exposition to how one might interpret these quantities. We start by noting that in some cases the norm of a density is linked to the precision parameter, as can be seen in the following example.

**Example 2.** Let  $U \sim \text{Unif}(a, b)$  and let  $\pi(x) = (b - a)^{-1}I_{(a,b)}(x)$  denote its corresponding density. Then, it holds that  $\|\pi\| = (\tau_U/12)^{1/4}$ , where the precision of  $U$  is  $\tau_U = 12/(b - a)^2$ . Next, consider a Normal model  $X \sim N(\mu, \tau_X)$  with known precision  $\tau_X$  and let  $\phi$  denote its corresponding density. It can be shown that  $\|\phi\| = \{\int_{\mathbb{R}} \phi^2(x; \mu, \tau_X) d\mu\}^{1/2} = \{\tau_X/(4\pi)\}^{1/4}$  which is a function of  $\tau_X$ .

The following proposition further explores the connection between norms and precision suggested by Example 2

**Proposition 1.** Let  $\Theta \subset \mathbb{R}^p$  with  $|\Theta| < \infty$  where  $|\cdot|$  denotes the Lebesgue measure. Consider  $\pi : \Theta \rightarrow [0, \infty)$  a probability density with  $\pi \in L_2(\Theta)$  and let  $\pi_0 \sim \text{Unif}(\Theta)$  denote a Uniform density on  $\Theta$ , then

$$\|\pi\|^2 = \|\pi - \pi_0\|^2 + \|\pi_0\|^2. \quad (7)$$

Since  $\|\pi_0\|^2$  is constant,  $\|\pi\|^2$  increases as  $\pi$ 's mass becomes more concentrated (or less Uniform). Thus, as can be seen from (7),  $\|\pi\|$  is a measure of how much  $\pi$  differs from a Uniform distribution over  $\Theta$ . This interpretation cannot be applied to  $\Theta$ 's that are not finite measurable as there is no corresponding proper Uniform distribution. Nonetheless, the notion that the norm of a density is a measure of its peakedness may be applied whether or not  $\Theta$  is finite measurable. Therefore,  $\|\cdot\|$  can be seen as very simple alternative to that proposed in Evans and Jang (2011) to compare the ‘informativeness’ of two competing priors with  $\|\pi_1\| < \|\pi_2\|$  indicating that  $\pi_1$  is less informative.

Further reinforcing the idea that the norm is related to the peakedness of a distribution, there is an interesting connection between  $\|\pi\|$  and the (differential) entropy (denoted by  $H_\pi$ ) which is described in the following theorem.

**Theorem 1.** Suppose  $\pi \in L_2(\Theta)$  is a continuous density on a compact  $\Theta \subset \mathbb{R}^p$ , and that  $\pi(\boldsymbol{\theta})$  is differentiable on  $\text{int}(\Theta)$ . Let  $H_\pi = -\int_{\Theta} \pi(\boldsymbol{\theta}) \log \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}$ . Then, it holds that

$$\|\pi\|^2 = 1 - H_\pi + o\{\pi(\boldsymbol{\theta}^*) - 1\}, \quad (8)$$

for some  $\boldsymbol{\theta}^* \in \text{int}(\Theta)$ .



The expansion in (8) hints that the norm of a density and the entropy should be negatively related, and hence as the norm of a density increases, its mass becomes more concentrated. In terms of priors, this suggests that priors with a large norm should be more ‘peaked’ relative to priors with a smaller norm. Therefore, the magnitude of a prior appears to be linked to its peakedness (as is demonstrated in (7) and in Example 2). While this might also be viewed as ‘informativeness,’ the Beta( $a, b$ ) density has a higher norm if  $(a, b) \in (1/2, 1)^2$  than if  $a = b = 1$ , possibly placing this interpretation at odds with the notion that  $a$  and  $b$  represent ‘prior successes’ and ‘prior failures’ in the Beta-Binomial setting.

As can be seen from (8), the connection between entropy and  $\|\pi\|$  is an approximation at best. Just as a first order Taylor expansion provides a poor polynomial approximation for points that are far from the point under which the expansion is made, the expansion in (8) will provide a poor entropy approximation when  $\pi$  is not similar to a standard Uniform-like distribution  $\pi_0$ . However, since  $\|\pi_0\|^2 = 1 - H_{\pi_0}$ , the approximation is exact for a standard Uniform-like distribution. We end this discussion by noting that integrals related to  $\|\pi\|^2$  also appear in physical models on  $L_2$ -spaces and they are usually interpreted as the total energy of a physical system (Hunter and Nachtergaele, 2005, p. 142).

Now, to illustrate the information that  $\|\cdot\|$  and  $\kappa$  provide, we consider the example described in the Introduction.

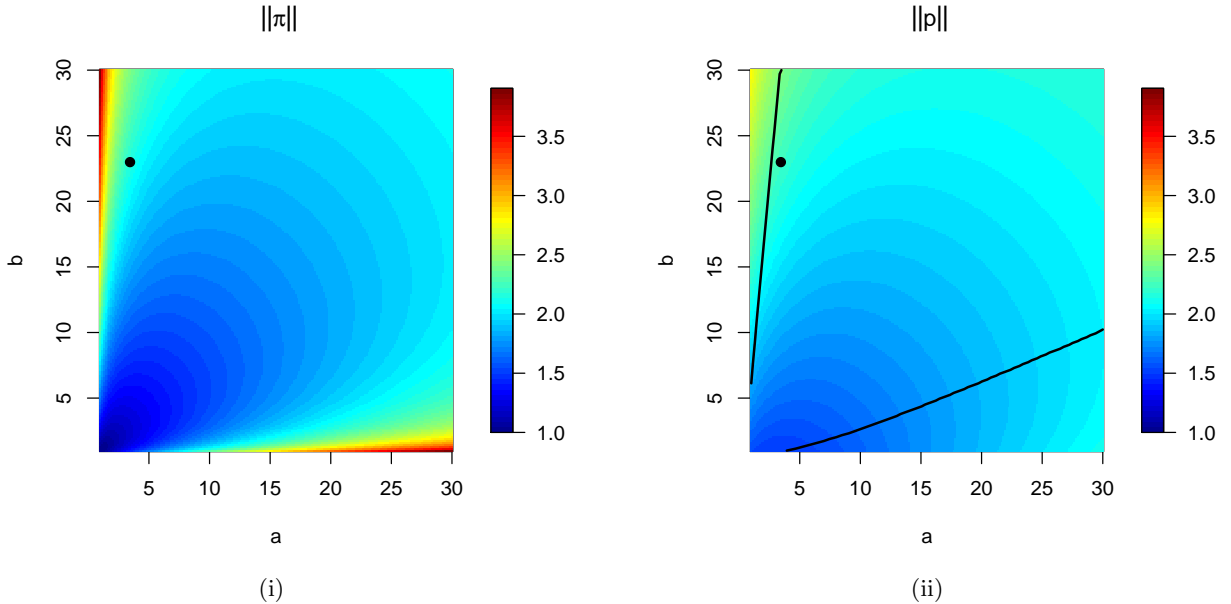


Figure 2: Prior and posterior norms for on-the-job drug usage toy example. Contour plots depicting the  $\|\cdot\|$  associated with a Beta( $a, b$ ) prior (i) and the corresponding Beta( $a^*, b^*$ ) posterior (ii), with  $a^* = a + 2$  and  $b^* = b + 8$ . Solid lines in (ii) indicate boundaries delimiting the region of values of  $a$  and  $b$  for which  $\|\pi\| > \|p\|$ . The solid dot ( $\bullet$ ) corresponds to  $(a, b) = (3.44, 22.99)$  (values employed by Christensen et al. 2011, pp. 26–27).

**Example 3** (On-the-job drug usage toy example, cont. 1). From the example in the Introduction we have  $\theta \mid \mathbf{y} \sim \text{Beta}(a^*, b^*)$  with  $a^* = a + n_1 = a + 2$  and  $b^* = b + n - n_1 = b + 8$ . The norm of the prior, posterior, and likelihood are respectively given by

$$\|\pi(a, b)\| = \frac{\{B(2a - 1, 2b - 1)\}^{1/2}}{B(a, b)}, \quad \|p(a, b)\| = \|\pi(a^*, b^*)\|,$$

with  $a, b > 1/2$ , and

$$\|\ell\| = \binom{n}{n_1} [B(2n_1 + 1, 2(n - n_1) + 1)]^{1/2},$$

where  $B(a, b) = \int_0^1 u^{a-1}(1-u)^{b-1} du$ .

Figure 2 (i) plots  $\|\pi(a, b)\|$  and Figure 2 (ii) plots  $\|p(a, b)\|$  as functions of  $a$  and  $b$ . We highlight the prior values  $(a_0, b_0) = (3.44, 22.99)$  which were employed by Christensen et al. (2011). Because prior densities with large norms will be more peaked relative to priors with small norms,  $\|\pi(a_0, b_0)\| = 2.17$  is more peaked than  $\|\pi(1, 1)\| = 1$  (Uniform prior) indicating that  $\|\pi(a_0, b_0)\|$  is more ‘informative’ than  $\|\pi(1, 1)\|$ . The norm of the posterior for these same pairs is  $\|p(a_0, b_0)\| = 2.24$  and  $\|p(1, 1)\| = 1.55$ , meaning that the posteriors will have mass more concentrated than the corresponding priors. In fact, the lines found in Figure 2 (ii) represent boundary lines such that all  $(a, b)$  pairs that fall outside of the boundary produce  $\|\pi(a, b)\| > \|p(a, b)\|$  which indicates that the prior is more peaked than the posterior (typically an undesirable result). If we used an extremely peaked prior, say  $(a_1, b_1) = (40, 300)$ , then we would get  $\|\pi(a_1, b_1)\| = 4.03$  and  $\|p(40, 300)\| = 4.04$  indicating that the peakedness of the prior and posterior densities is essentially the same.

Considering  $\kappa_{\pi, \ell}$ , it follows that

$$\kappa_{\pi, \ell}(a, b) = \frac{B(a^*, b^*)}{\{B(2a - 1, 2b - 1)B(2n_1 + 1, 2(n - n_1) + 1)\}^{1/2}}. \quad (9)$$

Figure 3 (i) plots values of  $\kappa$  as a function of prior parameters  $a$  and  $b$  with  $\kappa_{\pi, \ell}(a_0, b_0) \approx 0.69$  being highlighted indicating a great deal of agreement with the likelihood. In this example a lack of prior-data compatibility would occur (e.g.,  $\kappa_{\pi, \ell} \leq 0.1$ ) for priors that are very peaked at  $\theta > 0.95$  or for priors that place substantial mass at  $\theta < 0.5$ .

The values of the hyperparameters  $(a, b)$  which, according to  $\kappa_{\pi, \ell}$ , are more compatible with the data (i.e., those that maximize  $\kappa$ ) are given by  $(a^*, b^*) = (3, 9)$  and are highlighted with a star (\*) in Figure 3 (i). In Section 2.4 we provide some connections between this prior and maximum likelihood estimators.

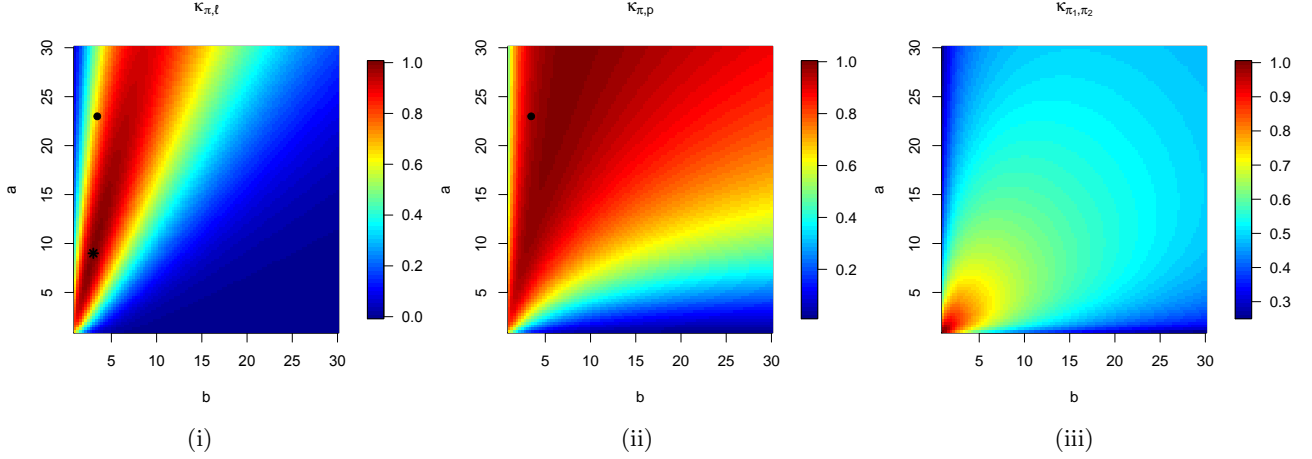


Figure 3: Compatibility ( $\kappa$ ) for on-the-job drug usage toy illustration as found in Equation (9) and Example 4. (i) Prior-likelihood compatibility,  $\kappa_{\pi,\ell}(a, b)$ ; the black star (\*) corresponds to  $(a^*, b^*)$  which maximize  $\kappa_{\pi,\ell}(a, b)$ . (ii) Prior-posterior compatibility,  $\kappa_{\pi,p}(a, b)$ . (iii) Prior-prior compatibility,  $\kappa_{\pi_1,\pi_2}(1, 1, a, b)$ , where  $\pi_1 \sim \text{Beta}(1, 1)$  and  $\pi_2 \sim \text{Beta}(a, b)$ . In (i) and (ii) the solid dot (•) corresponds to  $(a, b) = (3.44, 22.99)$  (values employed by Christensen et al. 2011, pp. 26–27).

### 2.3 Angles between other vectors

As mentioned, we are not restricted to use  $\kappa$  only to compare  $\pi$  and  $\ell$ . In fact, angles between different densities, and between likelihoods and densities or even between two likelihoods are available. We explore these options further using the example provided in the Introduction.

**Example 4** (On-the-job drug usage toy example, cont. 2). Extending Example 3 and Equation (9) we calculate

$$\kappa_{\pi,p}(a, b) = \frac{B(n_1 + 2a - 1, n - n_1 + 2b - 1)}{\{B(2a - 1, 2b - 1)B(2n_1 + 2a - 1, 2n - 2n_1 + 2b - 1)\}^{1/2}},$$

and for  $\pi_1 \sim \text{Beta}(a_1, b_1)$  and  $\pi_2 \sim \text{Beta}(a_2, b_2)$ ,

$$\kappa_{\pi_1,\pi_2}(a_1, b_1, a_2, b_2) = \frac{B(a_1 + a_2 - 1, b_1 + b_2 - 1)}{\{B(2a_1 - 1, 2b_1 - 1)B(2a_2 - 1, 2b_2 - 1)\}^{1/2}}.$$

To visualize how the hyperparameters influence  $\kappa_{\pi,p}$  and  $\kappa_{\pi_1,\pi_2}$  we provide Figures 3 (ii) and (iii). Figure 3 (ii) again highlights the prior used in Christensen et al. (2011) with  $\kappa_{\pi,p}(a_0, b_0) \approx 0.95$ ; see solid dot (•). This value of  $\kappa_{\pi,p}$  implies that both prior and posterior are concentrated on essentially the same subset of  $[0, 1]$ , indicating a large amount of agreement between them. Disagreement between prior and posterior takes place with priors concentrated on high probabilities of  $\theta$  being greater than 0.8. In Figure 3 (iii),  $\kappa_{\pi_1,\pi_2}$  is largest when  $\pi_2$  is close to  $\text{Unif}(0, 1)$  (the distribution of  $\pi_1$ ) and gradually drops off as  $\pi_2$  becomes more peaked and/or less symmetric.

In the next example, we utilize another small data set and demonstrate the application of  $\kappa$  to a two-parameter model.

**Example 5.** Hoff (2009, pp. 72–76) used a data set of nine midge wing lengths (originally reported by Grogan and Wirth 1981). The nine measurements were assumed to be conditionally iid with a Normal likelihood. The prior distribution for  $\mu$  and  $\sigma^2$  was decomposed as a Normal-Inverse Gamma distribution, i.e.,  $\mu \mid \sigma^2 \sim N(\mu_0, \sigma^2/\eta_0)$  and  $\sigma^2 \sim \text{IG}(\nu_0/2, \sigma_0^2\nu_0/2)$ ; we refer to this conjugate prior distribution as  $\text{NormIG}(\mu_0, \eta_0, \nu_0, \sigma_0^2)$ . As noted by Hoff (2009, p. 74), this parametrization affords appealing interpretations for the hyperparameters:  $\mu_0$  and  $\eta_0$  as the mean and sample size of ‘prior observations’—for inference on  $\mu \mid \sigma^2$ —and  $\nu_0$  and  $\sigma_0^2$  as the sample size and variance of ‘prior observations’—for inference on  $\sigma^2$ . In comparing two Normal-Inverse Gamma distributions,  $\text{NormIG}(\mu_1, \eta_1, \nu_1, \sigma_1^2)$  and  $\text{NormIG}(\mu_2, \eta_2, \nu_2, \sigma_2^2)$ ,  $\kappa_{\pi_1, \pi_2}$  may be expressed using the Normal-Inverse Gamma density with three different sets of hyperparameters, each evaluated at  $(\mu = 0, \sigma^2 = 1)$ , i.e.,

$$\kappa_{\pi_1, \pi_2} = \frac{(\pi_A \pi_B)^{1/2}}{\pi_C} \Big|_{\mu=0, \sigma^2=1}. \quad (10)$$

In this form,  $\pi_A$  represents the  $\text{NormIG}(\mu_1, 2\eta_1, 2\nu_1+3, \nu_1\sigma_1^2/(\nu_1+3/2))$  density,  $\pi_B$  the  $\text{NormIG}(\mu_2, 2\eta_2, 2\nu_2+3, \nu_2\sigma_2^2/(\nu_2+3/2))$  density, and  $\pi_C$  the  $\text{NormIG}((\eta_1\mu_1 + \eta_2\mu_2)/(\eta_1 + \eta_2), \eta_1 + \eta_2, \nu_1 + \nu_2 + 3, \{\nu_1\sigma_1^2 + \nu_2\sigma_2^2 + \eta_1\eta_2(\mu_1 - \mu_2)^2/(\eta_1 + \eta_2)\}/(\nu_1 + \nu_2 + 3))$  density.

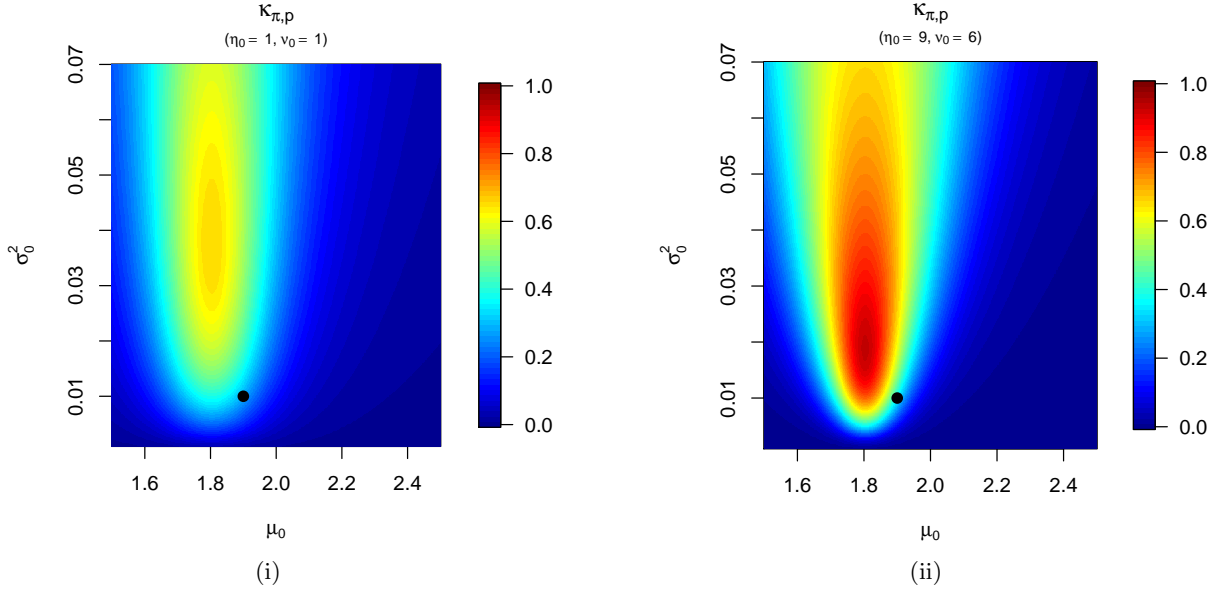


Figure 4: Prior-posterior compatibility,  $\kappa_{\pi,p}(\mu_0, \eta_0, \nu_0, \sigma_0^2)$ , for midge wing lengths data from Example 5. In (i)  $\eta_0$  and  $\nu_0$  are fixed at one, whereas in (ii)  $\eta_0$  is fixed at nine and  $\nu_0$  is fixed at six. The solid dot ( $\bullet$ ) corresponds to  $(\mu_0, \sigma_0^2) = (1.9, 0.01)$  which is here used as a baseline given that hyperparameters employed by Hoff (2009, pp. 72–76) are  $\mu_0 = 1.9, \eta_0 = 1, \nu_0 = 1$ , and  $\sigma_0^2 = 0.01$ .

In particular, (10) may be used not only for assessing agreement between two Normal-Inverse Gamma priors, but also between the prior and the posterior distribution. The hyperparameters for

the posterior relate to the prior specification as follows (see also Hoff, 2009, p. 75):

$$\begin{cases} \mu^* = (n\bar{Y} + \eta_0\mu_0)/(n + \eta_0), & \eta^* = \eta_0 + n, & \nu^* = \nu_0 + n, \\ \sigma^{2*} = \left\{ \nu_0\sigma_0^2 + \sum_{i=1}^n (Y_i - \bar{Y})^2 + \eta_0n(\eta^*)^{-1}(\mu_0 - \bar{Y})^2 \right\} / \nu^*. \end{cases}$$

For the midge data application, Hoff chose as hyperparameters  $\mu_0 = 1.9$ ,  $\eta_0 = 1$ ,  $\nu_0 = 1$ , and  $\sigma_0^2 = 0.01$ , while  $\bar{Y} = 1.804$  and  $\sum_{i=1}^n (Y_i - \bar{Y})^2 \approx 0.135$ , producing  $\kappa_{\pi,p} \approx 0.28$ . The agreement between the prior and posterior is not particularly strong. Figure 4 (i) displays the prior-posterior compatibility,  $\kappa_{\pi,p}$ , for these data as a function of  $\mu_0$  and  $\sigma_0^2$  while fixing  $\nu_0 = 1$  and  $\eta_0 = 1$ . To evaluate how  $\kappa_{\pi,p}$  is affected by  $\nu_0$  and  $\eta_0$ , the sample sizes of ‘prior observations,’ the analogous plot is displayed as Figure 4 (ii) when these values are fixed at  $\nu_0 = 6$  and  $\eta_0 = 9$ ; these alternative values for  $\nu_0$  and  $\eta_0$  are those which allow the compatibility between the prior and likelihood to be maximized. It is apparent from these plots that a somewhat larger value of  $\sigma_0^2$  would have increased  $\kappa_{\pi,p}$  substantially, and a simultaneous increase of  $\nu_0$  and  $\eta_0$  would further propel this increase.

## 2.4 Max-compatible priors and maximum likelihood estimators

In Example 3 we briefly alluded to a connection between priors maximizing prior-likelihood compatibility  $\kappa_{\pi,\ell}$ —to be termed as max-compatible priors—and maximum likelihood (ML) estimators, on which we now elaborate. In the following we use the notation  $\pi(\boldsymbol{\theta} \mid \boldsymbol{\alpha})$  to denote a prior on  $\boldsymbol{\theta} \in \Theta$ , and where  $\boldsymbol{\alpha} \in \mathcal{A}$  are hyperparameters. (Think of the Beta-Binomial model, where  $\theta \in \Theta = (0, 1)$ , and  $\boldsymbol{\alpha} = (a, b) \in \mathcal{A} = (0, \infty)^2$ .) Below, let  $\dim(\mathcal{A}) = q$  and  $\dim(\Theta) = p$ .

**Definition 3** (Max-compatible prior). *Let  $\mathbf{y} \sim f(\cdot \mid \boldsymbol{\theta})$ , and let  $\mathcal{P} = \{\pi(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) : \boldsymbol{\alpha} \in \mathcal{A}\}$  be a family of priors for  $\boldsymbol{\theta}$ . If there exists  $\boldsymbol{\alpha}_{\mathbf{y}}^* \in \mathcal{A}$ , such that  $\kappa_{\pi,\ell}(\boldsymbol{\alpha}_{\mathbf{y}}^*) = 1$ , the prior  $\pi(\boldsymbol{\theta} \mid \boldsymbol{\alpha}_{\mathbf{y}}^*) \in \mathcal{P}$  is said to be max-compatible, and  $\boldsymbol{\alpha}_{\mathbf{y}}^*$  is said to be a max-compatible hyperparameter.*

The max-compatible hyperparameter,  $\boldsymbol{\alpha}_{\mathbf{y}}^*$ , is by definition a random vector, and thus a max-compatible prior density is a random function. Geometrically, a prior is max-compatible iff it is collinear to the likelihood in the sense that  $\kappa_{\pi,\ell}(\boldsymbol{\alpha}_{\mathbf{y}}^*) = 1$  iff  $\pi(\boldsymbol{\theta} \mid \boldsymbol{\alpha}_{\mathbf{y}}^*) \propto f(\mathbf{y} \mid \boldsymbol{\theta})$ , for all  $\boldsymbol{\theta} \in \Theta$ .

The following example suggests there could be a connection between the ML estimator of  $\boldsymbol{\theta}$  and the max-compatibility parameter  $\boldsymbol{\alpha}_{\mathbf{y}}^*$ .

**Example 6** (Beta-Binomial). Let  $\sum_{i=1}^n Y_i \sim \text{Bin}(n, \theta)$ , and suppose  $\theta \sim \text{Beta}(a, b)$ . Let

$$\mathcal{P} = \{\beta(\theta \mid a, b) : (a, b) \in (1/2, \infty)^2\},$$

where  $\beta(\theta \mid a, b) = \theta^{a-1}(1-\theta)^{b-1}/B(a, b)$ . It can be shown that the max-compatible prior is  $\pi(\theta \mid a^*, b^*) = \beta(\theta \mid a^*, b^*)$ , where  $a^* = 1 + n_1$ , and  $b^* = 1 + n - n_1$ , so that

$$\hat{\theta}_n = \arg \max_{\theta \in (0,1)} f(\mathbf{y} \mid \theta) = \bar{Y} = \frac{a^* - 1}{a^* + b^* - 2} =: m(a^*, b^*). \quad (11)$$

A natural question is whether there always exists a function  $m : \mathcal{A} \rightarrow \Theta$ , linking the max-compatible parameter with the ML estimator, as in (11)? The following theorem addresses this question.

**Theorem 2.** *Let  $\mathbf{y} \sim f(\cdot \mid \theta)$ , and let  $\mathcal{P} = \{\pi(\theta \mid \alpha) : \alpha \in \mathcal{A}\}$  be a family of priors for  $\theta$ . Suppose there exists a max-compatible prior  $\pi(\theta \mid \alpha_{\mathbf{y}}^*) \in \mathcal{P}$ , which we assume to be unimodal. Then,*

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} f(\mathbf{y} \mid \theta) = m_{\pi}(\alpha_{\mathbf{y}}^*) := \arg \max_{\theta \in \Theta} \pi(\theta \mid \alpha_{\mathbf{y}}^*).$$

Theorem 2 states that the mode of the max-compatible prior coincides with the ML estimator. Note that in Example 6,  $m(a^*, b^*) = (a^* - 1)/(a^* + b^* - 2)$  is indeed the mode of a Beta prior. The next examples illustrate further this result.

**Example 7** (Exp–Gamma). In this case the max-compatible prior is given by  $f_{\Gamma}(\theta \mid a^*, b^*) = b^{*a^*} / \Gamma(a^*) \theta^{a^*-1} \exp\{-b^*\theta\} I_{(0,\infty)}(\theta)$ , where  $(a^*, b^*) = (1 + n, \sum_{i=1}^n Y_i)$ . The connection with the ML estimator is the following

$$\hat{\theta} = \arg \max_{\theta \in \Theta} f(\mathbf{y} \mid \theta) = \frac{n}{\sum_{i=1}^n Y_i} = \frac{a^* - 1}{b^*} =: m_2(a^*, b^*). \quad (12)$$

**Example 8** (Poisson–Gamma). In this case the max-compatible prior is given by  $f_{\Gamma}(\theta \mid a^*, b^*) = b^{*a^*} / \Gamma(a^*) \theta^{a^*-1} \exp\{-b^*\theta\} I_{(0,\infty)}(\theta)$ , where  $(a^*, b^*) = (1 + \sum_{i=1}^n Y_i, n)$ . The max-compatible hyperparameter in this case is different from the one in Example 7, but still a similar connection holds

$$\hat{\theta} = \arg \max_{\theta \in \Theta} f(\mathbf{y} \mid \theta) = \bar{Y} = \frac{a^* - 1}{b^*} =: m_2(a^*, b^*).$$

The preceding examples and theorem suggest a possible connection between max-compatible priors and empirical Bayes priors. It is true that both are data-driven prior distributions but save a few special cases the max-compatible prior does not coincide with an empirical Bayes prior. Theorem 2 highlights a situation when they do. If an ML estimator is employed to provide prior parameter values in an empirical Bayes prior and the prior distribution is symmetric so that the mean and median are equal, then the max-compatible prior and empirical Bayes prior will be the same (e.g., a  $N(\mu, \sigma^2)$  prior with  $\mu = \bar{Y}$  and  $\sigma^2$  known). Therefore, the max-compatible prior can be thought of as an alternative data-based prior construction.

### 3 Posterior schemes and Bayes geometries

By rewriting Bayes theorem as in (2), it is natural to pose the question: Do other inner products exist that can be used to mimic the geometric principles described in Section 2 and yet produce inference different but related to the Bayesian paradigm? As we shall see next, the answer to this question is positive, and we will refer to such approaches as posterior schemes. Below, let  $\mathcal{Y} = \{\mathbf{y} : f(\mathbf{y} | \boldsymbol{\theta}) > 0\}$ .

**Definition 4** (Posterior scheme). *Let  $h_{\pi, \ell} : \Theta \times \mathcal{Y} \rightarrow (0, \infty)$  be a mapping. Let  $\mathcal{H} = (\Theta, \langle\langle \cdot, \cdot \rangle\rangle)$  be an inner product space, such that  $\langle\langle \pi, \ell \rangle\rangle = \int_{\Theta} h_{\pi, \ell}(\boldsymbol{\theta}, \mathbf{y}) d\boldsymbol{\theta}$ . A posterior scheme is a mapping  $\mathbf{p} : \Theta \times \mathcal{Y} \rightarrow (0, \infty)$ , defined as*

$$\mathbf{p}(\boldsymbol{\theta} \parallel \mathbf{y}) := \frac{h_{\pi, \ell}(\boldsymbol{\theta}, \mathbf{y})}{\langle\langle \pi, \ell \rangle\rangle}. \quad (13)$$

The simplest posterior scheme is defined through Bayes theorem; by setting  $h_{\pi, \ell}(\boldsymbol{\theta}, \mathbf{y}) = \pi(\boldsymbol{\theta})f(\mathbf{y}; \boldsymbol{\theta}) = \pi(\boldsymbol{\theta})\ell(\boldsymbol{\theta})$ , and thus  $\langle\langle \pi, \ell \rangle\rangle = \langle \pi, \ell \rangle = \int_{\Theta} \pi(\boldsymbol{\theta})\ell(\boldsymbol{\theta}) d\boldsymbol{\theta}$ . Thus, we obtain  $\mathbf{p}(\boldsymbol{\theta} \parallel \mathbf{y}) = p(\boldsymbol{\theta} | \mathbf{y})$ , and hence the posterior scheme corresponding to Bayes theorem (i.e., corresponding to the standard  $L_2$  inner product) is simply the posterior density.

Note that a posterior scheme provides a well defined probability distribution for  $\boldsymbol{\theta}$  as it always integrates to one when we integrate over  $\Theta$ , i.e.,

$$\int_{\Theta} \mathbf{p}(\boldsymbol{\theta} \parallel \mathbf{y}) d\boldsymbol{\theta} = \int_{\Theta} \frac{h_{\pi, \ell}(\boldsymbol{\theta}, \mathbf{y})}{\langle\langle \pi, \ell \rangle\rangle} d\boldsymbol{\theta} = \frac{\langle\langle \pi, \ell \rangle\rangle}{\langle\langle \pi, \ell \rangle\rangle} = 1. \quad (14)$$

It is clear from Definition 4 that to construct posterior schemes, all that one needs is to plug into (13) inner products that can be expressed as integrals. In addition, the construction of Bayes-type estimators,  $\delta_{\pi}$  based on  $\mathbf{p}(\boldsymbol{\theta} \parallel \mathbf{y})$ , can be performed by minimizing the expected posterior scheme loss

$$\int_{\Theta} L(\boldsymbol{\theta}, \delta_{\pi}) \mathbf{p}(\boldsymbol{\theta} \parallel \mathbf{y}) d\boldsymbol{\theta},$$

where  $L : \Theta \times \mathcal{D} \rightarrow \mathbb{R}$ , is a loss function and  $\mathcal{D}$  is the space of all decision rules. Just as non-Euclidean distances find their application in geometry, we argue that posterior schemes other than the posterior density could provide other sensible ways to update the prior with data. One possible example is discussed next.

**Example 9** (Weighted posterior scheme). An alternative posterior scheme to Bayes theorem can be constructed by using the weighted inner product  $\langle \pi, f \rangle_w = \int_{\Theta} w(\boldsymbol{\theta})\pi(\boldsymbol{\theta})\ell(\boldsymbol{\theta}) d\boldsymbol{\theta}$ , where  $w(\boldsymbol{\theta}) > 0$  is a weighted function (Hunter and Nachtergaele, 2005, pp. 140–141), and it is given by

$$\mathbf{p}(\boldsymbol{\theta} \parallel \mathbf{y}) = \frac{w(\boldsymbol{\theta})\pi(\boldsymbol{\theta})\ell(\boldsymbol{\theta})}{\langle\langle \pi, \ell \rangle\rangle_w} \propto w(\boldsymbol{\theta})\pi(\boldsymbol{\theta})\ell(\boldsymbol{\theta}). \quad (15)$$

A particularly appealing interpretation for this scheme is as a model for combining expert opinion in prior elicitation. If  $w(\boldsymbol{\theta})$  represents a prior obtained from a second expert, then the posterior scheme in (15) provides a natural model for combining two independent priors  $\pi_1$  and  $\pi_2$ , i.e.,

$$\begin{aligned} p(\boldsymbol{\theta} \parallel \mathbf{y}) &= \frac{\pi_1(\boldsymbol{\theta})\pi_2(\boldsymbol{\theta})\ell(\boldsymbol{\theta})}{\int_{\Theta} \pi_1(\boldsymbol{\theta})\pi_2(\boldsymbol{\theta})\ell(\boldsymbol{\theta}) d\boldsymbol{\theta}} \\ &= \frac{\pi_1(\boldsymbol{\theta})\pi_2(\boldsymbol{\theta})\ell(\boldsymbol{\theta})}{\langle \pi_1, \ell \rangle \int_{\Theta} \pi_2(\boldsymbol{\theta})p(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}} \\ &= \frac{\pi_2(\boldsymbol{\theta})p_{\pi_1}(\boldsymbol{\theta} \mid \mathbf{y})}{\int_{\Theta} \pi_2(\boldsymbol{\theta})p_{\pi_1}(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}}, \end{aligned} \quad (16)$$

which is similar to Bayes theorem, but where the posterior based on  $\pi_1$  ( $p_{\pi_1}$ ) replaces the likelihood. To aid in the interpretation of (16), suppose that priors arrive sequentially, with  $\pi_1$  arriving firstly and  $\pi_2$  secondly. Thus in the second stage of the learning process associated with this posterior scheme it can be seen from (16) how the state of knowledge is updated. Note in addition that

$$p(\boldsymbol{\theta} \parallel \mathbf{y}) = \frac{\pi_1(\boldsymbol{\theta})p_{\pi_2}(\boldsymbol{\theta} \mid \mathbf{y})}{\int_{\Theta} \pi_1(\boldsymbol{\theta})p_{\pi_2}(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}} = \frac{\pi_2(\boldsymbol{\theta})p_{\pi_1}(\boldsymbol{\theta} \mid \mathbf{y})}{\int_{\Theta} \pi_2(\boldsymbol{\theta})p_{\pi_1}(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}},$$

where  $p_{\pi_2}$  is the posterior based on  $\pi_2$ , and hence the order of the learning based on this scheme is irrelevant. An alternative way of interpreting this scheme could be through Bayes theorem itself, which follows from observing that

$$p(\boldsymbol{\theta} \parallel \mathbf{y}) = \frac{\pi_1(\boldsymbol{\theta})\pi_2(\boldsymbol{\theta})\ell(\boldsymbol{\theta})}{\int_{\Theta} \pi_1(\boldsymbol{\theta})\pi_2(\boldsymbol{\theta})\ell(\boldsymbol{\theta}) d\boldsymbol{\theta}} = \frac{\pi^*(\boldsymbol{\theta})\ell(\boldsymbol{\theta})}{\int_{\Theta} \pi^*(\boldsymbol{\theta})\ell(\boldsymbol{\theta}) d\boldsymbol{\theta}} = p_{\pi^*}(\boldsymbol{\theta} \mid \mathbf{y}), \quad (17)$$

where  $\pi^*(\boldsymbol{\theta}) = \pi_1(\boldsymbol{\theta})\pi_2(\boldsymbol{\theta}) / \int_{\Theta} \pi_1(\boldsymbol{\theta})\pi_2(\boldsymbol{\theta}) d\boldsymbol{\theta}$ . In terms of sampling, the connection in (17) shows that updating the posterior scheme is equivalent to updating the posterior based on  $\pi^*$ .

More generally, if  $N$  independent sources of prior information,  $\{\pi_1, \dots, \pi_N\}$ , are available the posterior scheme is given as

$$p(\boldsymbol{\theta} \parallel \mathbf{y}) = \frac{\Pi(\boldsymbol{\theta})\ell(\boldsymbol{\theta})}{\langle \pi, \ell \rangle_{\Pi/\pi_1}}, \quad \boldsymbol{\theta} \in \Theta, \quad (18)$$

where  $\Pi(\boldsymbol{\theta}) = \prod_{j=1}^N \pi_j(\boldsymbol{\theta})$ . For example, if  $N$  researchers each have independent  $\text{Beta}(a_j, b_j)$  priors for  $\theta$  in the Bernoulli likelihood model, the posterior scheme in (18) is equivalent to the standard posterior scheme that utilizes a  $\text{Beta}(\sum_{j=1}^N a_j - (N-1), \sum_{j=1}^N b_j - (N-1))$  prior distribution. If  $a_j - 1$  and  $b_j - 1$  are interpreted as the  $j$ th researcher's equivalent prior observations of successes and failures, this prior reflects the aggregation of the disparate prior data.

Example 9 suggests that posterior schemes can be regarded as an alternative way to redirect the prior vector using data: The scheme in (15) obeys the strong likelihood principle, is a valid probability model (in the sense that it integrates to one), and obeys similar geometric principles to the ones discussed in Section 2. This leads us to the following concept.



**Definition 5** (Bayes geometry). *A Bayes geometry  $\mathcal{B}$  consists of an abstract geometry equipped with a posterior scheme. In addition, we say that  $\mathcal{B} = \{\mathcal{P}, \mathcal{L}, \mathbf{p}\}$  is the canonical Bayes geometry if the posterior scheme  $\mathbf{p}$  is simply a posterior density, i.e.  $\mathbf{p}(\boldsymbol{\theta} \parallel \mathbf{y}) = p(\boldsymbol{\theta} \mid \mathbf{y})$ .*

As mentioned above, to generate new posterior schemes, one simply needs to plug into (13) an inner product that can be written as an integral. Thus, beyond the canonical Bayes geometry and the Bayes geometry based on the weighted inner product, there are multitudinous possibilities for how one might use posterior schemes to recast the prior vector (update information) using the likelihood vector.

## 4 Posterior and prior mean-based estimators of $\kappa$ and $\|\cdot\|$

In many situations closed form estimators of  $\kappa$  and  $\|\cdot\|$  are not available. This leads to considering algorithmic techniques to obtain estimates. As most Bayes methods resort to using MCMC methods it would be appealing to express  $\kappa_{\cdot,\cdot}$  and  $\|\cdot\|$  as functions of posterior expectations and employ MCMC iterates to estimate them. For example,  $\kappa_{\pi,p}$  can be expressed as

$$\kappa_{\pi,p} = E_p \pi(\boldsymbol{\theta}) \left[ E_p \left\{ \frac{\pi(\boldsymbol{\theta})}{\ell(\boldsymbol{\theta})} \right\} E_p \{ \ell(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) \} \right]^{-1/2}, \quad (19)$$

where  $E_p(\cdot) = \int_{\Theta} \cdot p(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}$  is the expected value with respect to the posterior density. A natural Monte Carlo estimator would then be

$$\hat{\kappa}_{\pi,p} = \frac{1}{B} \sum_{b=1}^B \pi(\boldsymbol{\theta}^b) \left[ \left\{ \frac{1}{B} \sum_{b=1}^B \frac{\pi(\boldsymbol{\theta}^b)}{\ell(\boldsymbol{\theta}^b)} \right\} \left\{ \frac{1}{B} \sum_{b=1}^B \ell(\boldsymbol{\theta}^b) \pi(\boldsymbol{\theta}^b) \right\} \right]^{-1/2}, \quad (20)$$

where  $\boldsymbol{\theta}^b$  denotes the  $b$ th MCMC iterate of  $p(\boldsymbol{\theta} \mid \mathbf{y})$ . Consistency of such an estimator follows trivially by the ergodic theorem and the continuous mapping theorem, but there is an important issue regarding its stability. Unfortunately, (19) includes an expectation that contains  $\ell(\boldsymbol{\theta})$  in the denominator and therefore (20) inherits the undesirable properties of the so-called harmonic mean estimator (Newton and Raftery, 1994). It has been shown that even for simple models this estimator may have infinite variance (Raftery et al. 2007), and has been harshly criticized for, among other things, converging extremely slowly. Indeed, as argued by Wolpert and Schmidler (2012, p. 655):

“the reduction of Monte Carlo sampling error by a factor of two requires increasing the Monte Carlo sample size by a factor of  $2^{1/\varepsilon}$ , or in excess of  $2.5 \cdot 10^{30}$  when  $\varepsilon = 0.01$ , rendering [the harmonic mean estimator] entirely untenable.”

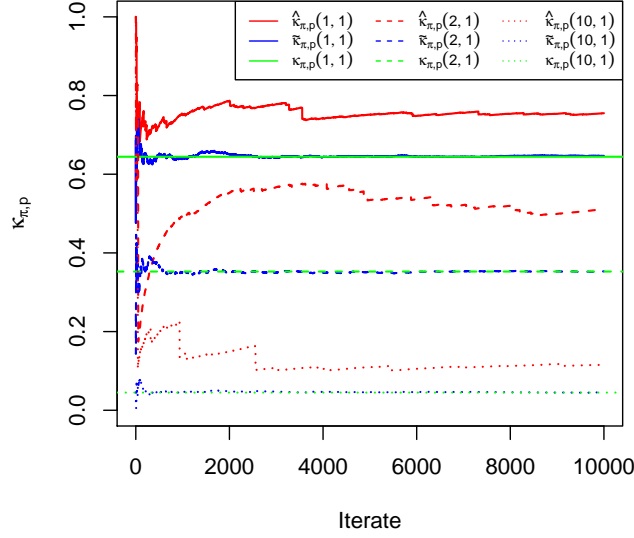


Figure 5: Running point estimates of prior-posterior compatibility,  $\kappa_{\pi,p}$ , for the on-the-job drug usage toy example. Green lines correspond to the true  $\kappa_{\pi,p}$  values computed as in Example 4, blue represents  $\tilde{\kappa}_{\pi,p}$  and red denotes  $\hat{\kappa}_{\pi,p}$ . Notice that  $\tilde{\kappa}_{\pi,p}$  converges to the true  $\kappa_{\pi,p}$  values quickly while  $\hat{\kappa}_{\pi,p}$  will need much more than 10 000 Monte Carlo draws to converge.

Making things a bit more difficult is the fact that (20) contains a *root* of  $1/\ell(\boldsymbol{\theta})$ , which renders corrections like those found in Lenk (2009) and Pajor and Osiewalski (2013) unsuitable.

An alternate strategy is to avoid writing  $\kappa_{\pi,p}$  as a function of harmonic mean estimators and instead express it as a function of posterior and prior expectations. For example, consider

$$\kappa_{\pi,p} = E_p \pi(\boldsymbol{\theta}) \left[ \frac{E_\pi \{\pi(\boldsymbol{\theta})\}}{E_\pi \{\ell(\boldsymbol{\theta})\}} E_p \{\ell(\boldsymbol{\theta}) \pi(\boldsymbol{\theta})\} \right]^{-1/2}, \quad (21)$$

where  $E_\pi(\cdot) = \int_{\Theta} \cdot \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}$ . Now the Monte Carlo estimator is

$$\tilde{\kappa}_{\pi,p} = \frac{1}{B} \sum_{b=1}^B \pi(\boldsymbol{\theta}^b) \left[ \left\{ \frac{\sum_{b=1}^B \pi(\boldsymbol{\theta}^b)}{\sum_{b=1}^B \ell(\boldsymbol{\theta}^b)} \right\} \left\{ \frac{1}{B} \sum_{b=1}^B \ell(\boldsymbol{\theta}^b) \pi(\boldsymbol{\theta}^b) \right\} \right]^{-1/2}, \quad (22)$$

where  $\boldsymbol{\theta}_b$  denotes the  $b$ th draw of  $\boldsymbol{\theta}$  from  $\pi(\boldsymbol{\theta})$ , which can also be sampled within the MCMC algorithm. Representations (21) and (22) are somewhat less elegant than (19) and (20) as they require draws from the posterior and the prior, but they behave much better in practice. To see this, Figure 5 contains running estimates of  $\kappa_{\pi,p}$  using (20) and (22) for Example 3 with three prior parameter specifications, namely:  $(a = 1, b = 1)$ ,  $(a = 2, b = 1)$ , and  $(a = 10, b = 1)$ ; the true  $\kappa_{\pi,p}$  for each prior specification is also provided. It is fairly clear that  $\hat{\kappa}_{\pi,p}$  displays slow convergence and large variance, while  $\tilde{\kappa}_{\pi,p}$  converges quickly.

The next proposition contains prior and posterior mean-based representations of geometric quantities for the canonical Bayes geometry, that can be readily used for constructing Monte Carlo estimators. Notice that metrics that include the prior only (e.g.,  $\|\pi\|$ ) are expressed entirely as functions of prior expectations. This allows comparing competing prior densities prior to any model fitting. (We briefly note that there is an enormous frequentist literature on the estimation of the integral  $\int_{\Theta} \pi^2(\boldsymbol{\theta}) d\boldsymbol{\theta}$ , especially due to its appearance in some variance–covariance structures; see for instance [Giné and Nickl \(2008\)](#), and the references therein.)

**Proposition 2.** *Let  $\mathcal{B} = \{\mathcal{P}, \mathcal{L}, \mathbf{p}\}$  be the canonical Bayes geometry and let  $\kappa$  denote compatibility in this geometry. Let  $E_p(\cdot) = \int_{\Theta} \cdot p(\boldsymbol{\theta} \mid \mathbf{y}) d\boldsymbol{\theta}$  and  $E_{\pi}(\cdot) = \int_{\Theta} \cdot \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}$  be the posterior and prior means. The following equalities hold:*

$$\begin{aligned} \|p\|^2 &= \frac{E_p\{\ell(\boldsymbol{\theta})\pi(\boldsymbol{\theta})\}}{E_{\pi}\ell(\boldsymbol{\theta})}, \quad \|\pi\|^2 = E_{\pi}\pi(\boldsymbol{\theta}), \quad \|\ell\|^2 = E_{\pi}\ell(\boldsymbol{\theta}) E_p\left\{\frac{\ell(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta})}\right\}, \\ \kappa_{\pi_1, \pi_2} &= E_{\pi_1}\pi_2(\boldsymbol{\theta}) \left[ E_{\pi_1}\pi_1(\boldsymbol{\theta}) E_{\pi_2}\pi_2(\boldsymbol{\theta}) \right]^{-1/2}, \quad \kappa_{\pi, \ell} = E_{\pi}\ell(\boldsymbol{\theta}) \left[ E_{\pi}\pi(\boldsymbol{\theta}) E_{\pi}\ell(\boldsymbol{\theta}) E_p\left\{\frac{\ell(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta})}\right\} \right]^{-1/2}, \\ \kappa_{\pi, p} &= E_p\pi(\boldsymbol{\theta}) \left[ \frac{E_{\pi}\pi(\boldsymbol{\theta})}{E_{\pi}\ell(\boldsymbol{\theta})} E_p\{\ell(\boldsymbol{\theta})\pi(\boldsymbol{\theta})\} \right]^{-1/2}, \quad \kappa_{\ell, p} = E_p\ell(\boldsymbol{\theta}) \left[ E_p\left\{\frac{\ell(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta})}\right\} E_p\{\ell(\boldsymbol{\theta})\pi(\boldsymbol{\theta})\} \right]^{-1/2}, \\ \kappa_{\ell_1, \ell_2} &= E_{\pi}\ell_2(\boldsymbol{\theta}) E_{p_2}\left\{\frac{\ell_1(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta})}\right\} \left[ E_{\pi}\{\ell_1(\boldsymbol{\theta})\} E_{p_1}\left\{\frac{\ell_1(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta})}\right\} E_{\pi}\ell_2(\boldsymbol{\theta}) E_{p_2}\left\{\frac{\ell_2(\boldsymbol{\theta})}{\pi(\boldsymbol{\theta})}\right\} \right]^{-1/2}. \end{aligned}$$

In the next section we provide an example that requires the use of Proposition 2 to estimate  $\kappa$  and  $\|\cdot\|$ .

## 5 Example: Regression shrinkage priors

### 5.1 Compatibility of Gaussian and Laplace priors

The linear regression model is ubiquitous in applied statistics. In vector form, the model is commonly written as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}). \quad (23)$$

where  $\mathbf{y} = (Y_1, \dots, Y_n)^T$ ,  $\mathbf{X}$  is a  $n \times p$  design matrix,  $\boldsymbol{\beta}$  is a  $p$ -vector of regression coefficients, and  $\sigma^2$  is an unknown idiosyncratic variance parameter. We consider two competing prior distributions for  $\boldsymbol{\beta}$ , Gaussian and Laplace. These two priors are often employed as shrinkage priors that perform some type of regularization. Connections between the regularization via ridge and lasso penalization and that from using Gaussian and Laplace prior distributions are now well documented ([Park and Casella 2008](#), [Kyung et al. 2010](#)). Estimating ridge regression coefficients amounts to minimizing  $\sum_{i=1}^n (Y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$

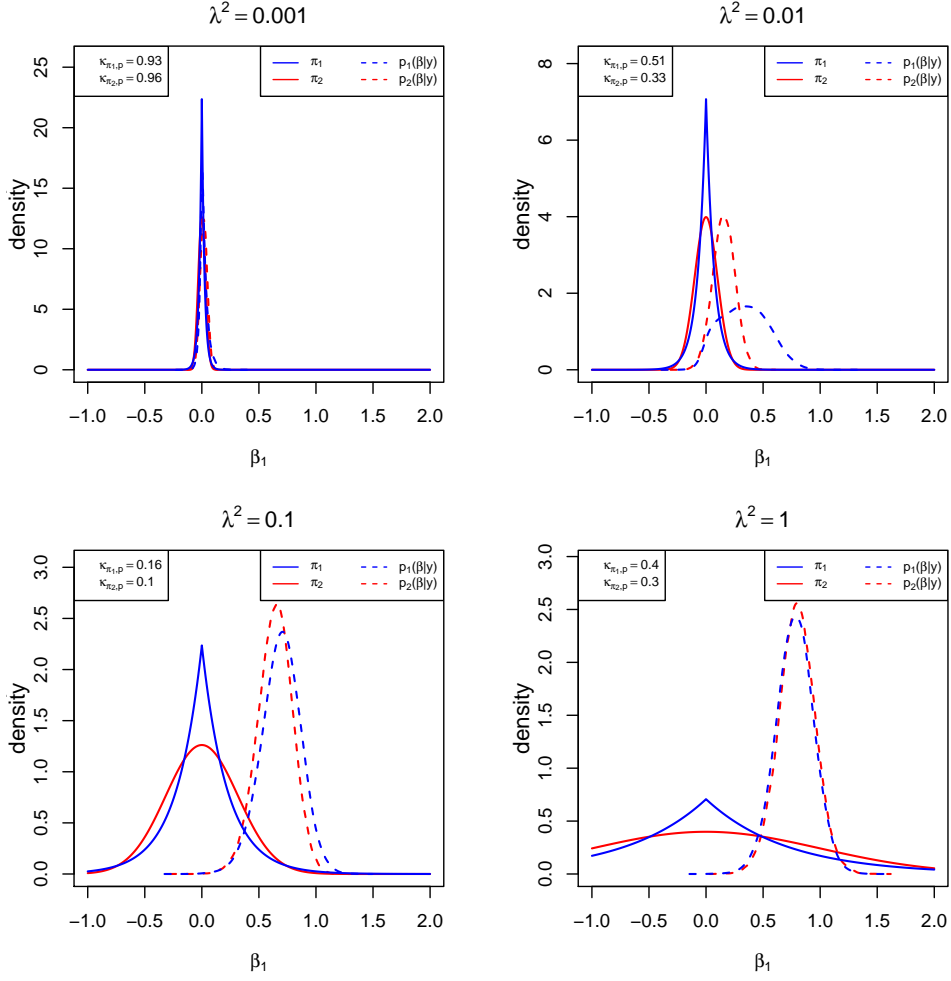


Figure 6: Prior and posterior distributions associated with  $\pi_1 \sim N(0, \lambda^2)$  and  $\pi_2 \sim \text{Laplace}(0, b)$  with  $b = \sqrt{0.5\lambda^2}$  ensuring that  $\text{var}_{\pi_1}(\beta_1) = \text{var}_{\pi_2}(\beta_1)$ . The plots vary in their  $\lambda$  values with  $\lambda^2 \in \{0.001, 0.01, 0.1, 1\}$ . The  $\kappa_{\pi, p}$  values provide an indication of mass intersection between prior and posterior.

subject to  $\sum_{j=1}^p \beta_j^2 < \lambda$  and it has been shown that assigning  $\beta_j \stackrel{\text{iid}}{\sim} N(0, \lambda^2)$  produces the same regularization on  $\beta$ . Similarly, estimating lasso coefficients amounts to minimizing  $\sum_{i=1}^n (Y_i - \mathbf{x}_i^T \beta)^2$  subject to  $\sum_{j=1}^p |\beta_j| < b$  and produces the same regularization as assigning  $\beta_j \stackrel{\text{iid}}{\sim} \text{Laplace}(0, b)$  with  $\text{var}(\beta_j) = 2b^2$ . In what follows we will use  $\pi_1$  to denote a Gaussian prior and  $\pi_2$  a Laplace. Further, to make reasonable comparisons between the two priors, we set  $b = \sqrt{0.5\lambda^2}$  which ensures that  $\text{var}_{\pi_1}(\beta_j) = \text{var}_{\pi_2}(\beta_j) = \lambda^2$  for all  $j$ . To develop intuition regarding these two priors we first present some results from a small synthetic dataset generated using a simple linear regression model. Then we consider the prostate cancer data example found in [Hastie, Tibshirani and Friedman \(2008, chap. 3.4\)](#) that was used to illustrate differences between the ridge and lasso regularization.

Using  $Y_i = X_i \beta_1 + \varepsilon_i$  with  $\beta_1 = 1$ ,  $\varepsilon_i \sim N(0, 2^2)$ , and  $X_i \sim \text{Unif}(-3, 3)$  as a data generating

mechanism we generated 25  $(X_i, Y_i)$  pairs. Using these 25 observations we fit the no-intercept simple linear regression model using  $\sigma \sim \text{UN}(0, 10)$  as a prior in addition to considering both  $\beta_1 \sim \pi_1$  and  $\beta_1 \sim \pi_2$ . Then, for each value of  $\lambda^2 \in \{0.001, 0.01, 0.1, 1\}$  we collected 10 000 MCMC draws from posteriors associated with  $\pi_1$  and  $\pi_2$  and employed them to compute  $\kappa_{\pi_1, p_1}$  and  $\kappa_{\pi_2, p_2}$  using Proposition 2. Results can be found in Figure 6 where the top left plot represents an extremely peaked prior for both  $\pi_1$  and  $\pi_2$  which produces posterior distributions that are very similar to the respective priors resulting in  $\kappa_{\pi, p}$  values that are close to one. As the value of  $\lambda^2$  increases, the priors become less ‘informative’ and agreement between prior and posterior decreases. However, in becoming more flat, the priors reach a point where the intersection of prior and posterior mass increases which is depicted by an increase in  $\kappa_{\pi, p}$  in the bottom right plot in Figure 6.

## 5.2 Prostate cancer data example

We now turn our attention to the prostate cancer data example found in [Hastie, Tibshirani and Friedman \(2008, chap. 3.4\)](#). In this example the response variable is the level of prostate-specific antigens measured on 97 males. Eight other clinical measurements (such as age and log prostate weight) were also measured and are used as covariates. Thus,  $p = 8$  in this example .

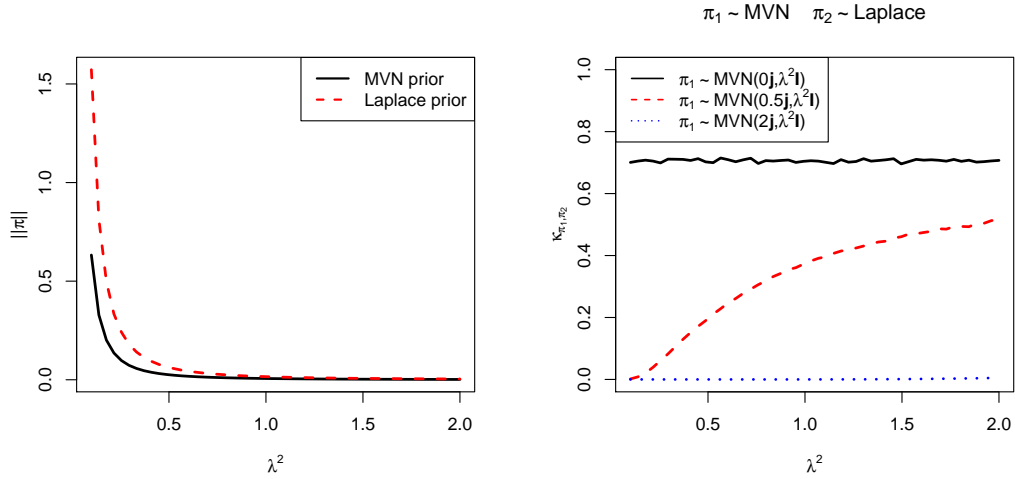


Figure 7: A comparison of priors associated with Ridge (MVN,  $\pi_1$ ) and Lasso (Laplace,  $\pi_2$ ) regularization in regression models in terms of  $\|\pi\|$  and  $\kappa_{\pi_1, \pi_2}$ . The left plot depicts  $\|\cdot\|$  as a function of  $\lambda^2$  for both  $\pi_1$  and  $\pi_2$ . The right compares  $\kappa_{\pi_1, \pi_2}$  values as a function of  $\lambda^2$  when  $\pi_1$  and  $\pi_2$  are centered at zero to that when the center of  $\pi_1$  moves away from zero.

Before proceeding with model fit, we first evaluate the ‘informativeness’ of the two priors for the eight regression coefficients by computing  $\|\pi_1\|$  and  $\|\pi_2\|$  and then assess their compatibility (or mass

intersection) by computing  $\kappa_{\pi_1, \pi_2}$ . All calculations employ Proposition 2. Each metric is calculated for a sequence of  $\lambda^2$  values with results provided in Figure 7. The left plot of Figure 7 provides  $\|\cdot\|$  of each prior for a sequence of  $\lambda^2$  values. For small values of the scale parameter  $\|\pi_1\| < \|\pi_2\|$ , indicating that the Laplace prior is more peaked than the Normal. Thus, even though the Laplace has thicker tails, it is more ‘informative’ relative to the Gaussian. This corroborates the lasso penalization’s ability to shrink coefficients to zero (something ridge regulation lacks). As  $\lambda^2$  increases the two norms converge as both spread their mass more uniformly. The right plot of Figure 7 depicts  $\kappa_{\pi_1, \pi_2}$  as a function of  $\lambda^2$ . When  $\pi_1$  is centered at zero, then  $\kappa_{\pi_1, \pi_2}$  is constant over values of  $\lambda^2$  which means that mass intersection when both priors are centered at zero is not influenced by tail thickness. Compare this to  $\kappa$  values when  $\pi_1$  is not centered at zero [i.e.,  $\pi_1 \sim \text{MVN}(0.5\mathbf{j}, \lambda^2\mathbf{I})$  or  $\pi_1 \sim \text{MVN}(2\mathbf{j}, \lambda^2\mathbf{I})$ ]. For the former,  $\kappa$  increases as intersection of prior and posterior mass increases. For the latter,  $\lambda^2$  must be greater than two for there to be any substantial mass intersection as  $\kappa_{\pi_1, \pi_2}$  remains essentially at zero.

Now that the ‘informativeness’ of the two priors has been explored, we fit model (23) to the cancer data. Within the MCMC algorithm we compute  $\kappa_{\pi_1, \ell}$  and  $\kappa_{\pi_2, \ell}$  along with  $\kappa_{\pi_1, p_1}$  and  $\kappa_{\pi_2, p_2}$  using Proposition 2. Without loss of generality we centered the  $\mathbf{y}$  so that  $\boldsymbol{\beta}$  does not include an intercept and standardized each of the eight covariates to have mean zero and standard deviation one. We employ  $\sigma \sim \text{UN}(0, 2)$  as a prior. The resulting  $\kappa$  for a range of  $\lambda^2$  values is provided in Figure 8. From Figure 8 it appears that prior-data agreement is very small for both priors indicating

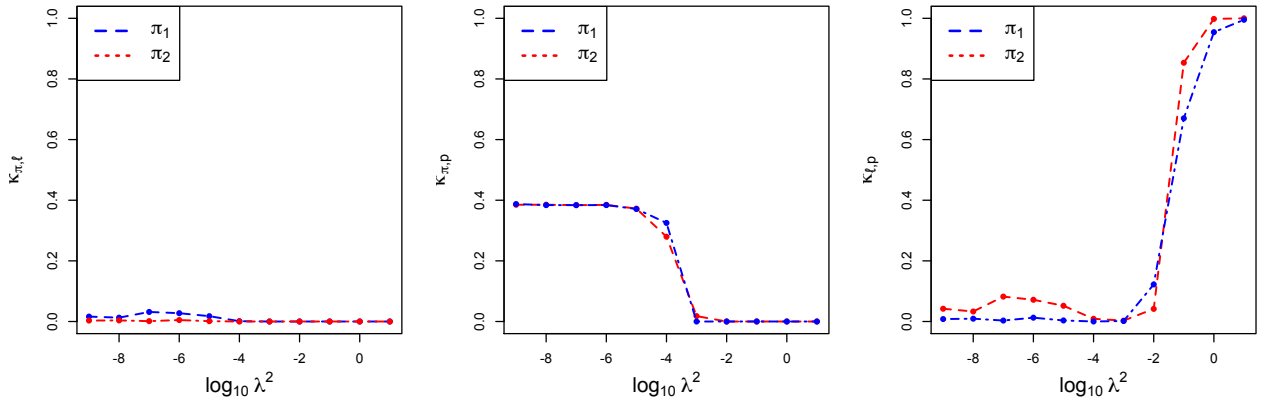


Figure 8: Compatibility ( $\kappa$ ) for linear regression model in (23), with shrinkage priors, applied to the prostate cancer data from Hastie, Tibshirani and Friedman (2008, chap. 3.4). The  $\kappa$  estimates were computed using Proposition 2.

the existence of prior-data incompatibility. However, for small values of  $\lambda^2$ ,  $\kappa_{\pi_1, \ell} > \kappa_{\pi_2, \ell}$  indicating

more compatibility between prior and data for the Gaussian prior. As an aside, the procedure found in [Evans and Moshonov \(2006\)](#) would conclude that no prior-data conflict exists in this example as the prior for  $\sigma$  is sufficiently diffuse to render the data plausible per the prior predictive distribution. Prior-posterior compatibility ( $\kappa_{\pi,p}$ ) is very similar for both priors with that for  $\pi_2$  being slightly smaller when  $\lambda^2$  is close to  $10^{-4}$ . Interestingly,  $\kappa_{\pi,p}$  for both priors appears to asymptote at around 0.4 as  $\lambda^2$  decreases. The value of the asymptote is an artifact of the prior on  $\sigma$ . The slightly higher  $\kappa_{\pi,p}$  value for the Gaussian prior implies that it has slightly more influence on the posterior than the Laplace, and  $\kappa_{p,\ell}$  communicates a similar story, mainly that the Gaussian prior has more influence on the posterior than the Laplace. The thicker tails of the Laplace prior seem to produce larger  $\kappa_{p,\ell}$  values than that of the Gaussian prior indicating a larger amount of posterior-data compatibility. Additionally,  $\kappa_{p_2,\ell}$  approaches one quicker than  $\kappa_{p_1,\ell}$ . This may be a result of the ability that the Laplace prior has to shrink coefficients to zero. Overall, it appears that the Gaussian prior has more influence on the resulting posterior distribution relative to the Laplace when updating knowledge via Bayes theorem.

## 6 Discussion

We discussed a natural geometric framework to Bayesian inference which motivated a simple, intuitively appealing measure of the agreement between priors, likelihoods, and posteriors: compatibility ( $\kappa$ ). In this geometric framework, we also discuss a related measure of the ‘informativeness’ of a distribution,  $\|\cdot\|$ . In addition, in [Section 4](#) we developed MCMC-based estimators of these metrics that are easily computable and, by avoiding the estimation of harmonic means, are reasonably stable. Therefore at virtually no cost, practitioners can easily produce metrics that assess the degree of prior-data and prior-posterior compatibility. Overall, we believe that the procedures developed in this paper should be a valuable contribution to the applied Bayesian modeling community.

In theory, one may argue that compatibility as defined in [Section 2](#) is grounded in the same construction principles as Pearson correlation, in the sense that both consist of standardized inner products. However, compatibility is defined for priors, posteriors, and likelihoods in  $L_2(\Theta)$  equipped with the inner product [\(1\)](#), whereas Pearson correlation works with random variables in  $L_2(\Omega, \mathbb{B}_\Omega, P)$  equipped with the inner product [\(6\)](#). Our concept of compatibility can be used to evaluate how much the prior agrees with the likelihood, to measure the sensitivity of the posterior to the prior, and to quantify the level of agreement of elicited priors. One practical drawback with our geometric construction is that it has been developed for priors which are on  $L_2(\Theta)$ , and thus some cases will not be handled by our setting; a simple example is that of the Jeffreys prior for the Beta-Binomial,

Beta(1/2, 1/2), whose norm will be infinity. One possible approach to being able to consider densities not in  $L_2(\Theta)$ , but not explored further here, is to work directly with  $\langle \sqrt{g}, \sqrt{h} \rangle = \int_{\Theta} \sqrt{g(\boldsymbol{\theta})} \sqrt{h(\boldsymbol{\theta})} d\boldsymbol{\theta}$ , for  $g, h \in L_1(\Theta)$ . This approach would result in a compatibility measure,  $\kappa_{\sqrt{g}, \sqrt{h}}$ , that continues being a metric that measures agreement between two elements of a geometry, though it loses direct connection with Bayes theorem. Interestingly,  $\kappa_{\sqrt{g}, \sqrt{h}}$  coincides with the so-called Hellinger affinity (van der Vaart, 1998, p. 211)

In addition to assessing agreement between the three components of Bayes theorem,  $\kappa$  can also be used to perform model comparison by comparing competing likelihoods. It can also be employed to assess sensitivity to potential influential points by comparing same likelihoods and/or posteriors with and without the data points under consideration.

The question of how one might gauge the compatibility between prior and likelihood in hierarchical models is a natural one. For example, to complete the hierarchy in Section 5 one could consider assigning a prior to  $\lambda^2$ . Without fully specifying the prior distributions on parameters that appear in the likelihood (i.e., introducing a process model in a hierarchy),  $\kappa_{\pi, \ell}$  would become a direct function of the process model parameters. To assess compatibility between prior and likelihood, it would then be natural to consider  $\int \kappa_{\pi, \ell}(\lambda^2) \pi(\lambda^2 | \mathbf{y}) d\lambda^2$ , where

$$\kappa_{\pi, \ell}(\lambda^2) = \frac{\int_{\Theta} \pi(\boldsymbol{\theta} | \lambda^2) \ell(\boldsymbol{\theta}) d\boldsymbol{\theta}}{\{\int_{\Theta} \pi^2(\boldsymbol{\theta} | \lambda^2) d\boldsymbol{\theta}\}^{1/2} \{\int_{\Theta} \ell^2(\boldsymbol{\theta}) d\boldsymbol{\theta}\}^{1/2}}.$$

This, however, would increase the computational cost of estimating  $\kappa$  considerably. An alternative approach might be to employ a plug-in summary of  $\kappa$  via  $\kappa_{\pi, \ell}(\lambda_{\mathbf{y}}^2)$  where  $\lambda_{\mathbf{y}}^2$  corresponds to either the posterior mean or the max-compatible hyperparameter. Assessing the value of this approach and considering  $\kappa$  in hierarchical models is the topic of current research.

A possible avenue of research which is only briefly explored here is that of using our geometry for devising new probabilistic models for recasting the prior vector using the likelihood vector. Indeed, as mentioned in Section 3, new posterior schemes may be generated by plugging into (13) alternative inner products which can be written as an integral. Thus, the sky really is the limit in how one might use posterior schemes to recast the prior vector using the likelihood vector. Developing innovative ways of recasting the prior vector is an area of ongoing research. Another possibility for future research is an exploration of how the dimensionality of  $\Theta$  should affect the interpretation of  $\kappa$ , if at all. We would anticipate that as the dimensionality increases, there is increased potential for disagreement between two distributions. Consequently,  $\kappa$  would generally diminish as additional parameters are added, *ceteris paribus*. A suitable offsetting transformation of  $\kappa$ , if it exists, could result in a measure of ‘per parameter’ agreement.



## Appendix

*Proof of Theorem 1.* The proof follows by combining a Taylor expansion with the first mean value theorem for integrals (Bartle and Sherbert, 2010, p. 301). Just note that

$$\begin{aligned}
H_\pi &= - \int_{\Theta} \pi(\boldsymbol{\theta}) \log \pi(\boldsymbol{\theta}) \, d\boldsymbol{\theta} \\
&= - \int_{\Theta} \pi(\boldsymbol{\theta}) [\pi(\boldsymbol{\theta}) - 1 + o\{\pi(\boldsymbol{\theta}) - 1\}] \, d\boldsymbol{\theta} \\
&= - \int_{\Theta} \pi^2(\boldsymbol{\theta}) \, d\boldsymbol{\theta} + \int_{\Theta} \pi(\boldsymbol{\theta}) \, d\boldsymbol{\theta} \\
&\quad - \int_{\Theta} \pi(\boldsymbol{\theta}) o\{\pi(\boldsymbol{\theta}) - 1\} \, d\boldsymbol{\theta} \\
&= 1 - \|\pi\|^2 + o\{\pi(\boldsymbol{\theta}^*) - 1\} \int_{\Theta} \pi(\boldsymbol{\theta}) \, d\boldsymbol{\theta} \\
&= 1 - \|\pi\|^2 + o\{\pi(\boldsymbol{\theta}^*) - 1\},
\end{aligned}$$

for some  $\boldsymbol{\theta}^* \in \text{int}(\Theta)$ , from where the final result follows.  $\square$

*Proof of Proposition 1.* Note that

$$\begin{aligned}
\|\pi\|^2 &= \int_{\Theta} \pi^2(\boldsymbol{\theta}) \, d\boldsymbol{\theta} \\
&= \int_{\Theta} \{\pi(\boldsymbol{\theta}) - \pi_0(\boldsymbol{\theta})\}^2 \, d\boldsymbol{\theta} + \int_{\Theta} \pi_0^2(\boldsymbol{\theta}) \, d\boldsymbol{\theta} \\
&= \|\pi - \pi_0\|^2 + \|\pi_0\|^2.
\end{aligned}$$

$\square$

*Proof of Theorem 2.* Just note that  $\kappa_{\pi, f}(\boldsymbol{\alpha}_y^*) = 1$  can be equivalently restated as  $\pi(\boldsymbol{\theta} \mid \boldsymbol{\alpha}_y^*) \propto f(\mathbf{y} \mid \boldsymbol{\theta})$ , for all  $\boldsymbol{\theta} \in \Theta$ , which in turn implies that  $\arg \max_{\boldsymbol{\theta} \in \Theta} \pi(\boldsymbol{\theta} \mid \boldsymbol{\alpha}_y^*) = \arg \max_{\boldsymbol{\theta} \in \Theta} f(\mathbf{y} \mid \boldsymbol{\theta})$ .  $\square$

*Proof of Proposition 2.* Since  $\langle \pi, \ell \rangle = E_\pi \ell(\boldsymbol{\theta})$ , it follows that

$$\|p\|^2 = \int_{\Theta} p^2(\boldsymbol{\theta} \mid \mathbf{y}) \, d\boldsymbol{\theta} = \frac{1}{\langle \pi, \ell \rangle} \int_{\Theta} \ell(\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathbf{y}) \, d\boldsymbol{\theta} = \frac{E_p \{\ell(\boldsymbol{\theta}) \pi(\boldsymbol{\theta})\}}{E_\pi \ell(\boldsymbol{\theta})}.$$

Similarly, we can derive  $\|\ell\|$ , so that the remaining results follow from the definition of  $\kappa$  in (5).  $\square$

## References

- Agarawal, A. and Daumé, III, H. (2010) A geometric view of conjugate priors. *Mach. Learn.*, **81**, 99–113.
- Aitchison, J. (1971) A geometrical version of Bayes’ theorem. *Am. Statist.*, **25**, 45–46.
- Al Labadi, L. and Evans, M. (2016) Optimal robustness results for relative belief inferences and the relationship to prior-data conflict. *Bayesian Anal.*, in press.
- Bartle, R. and Sherbert, D. (2010) *Introduction to Real Analysis* (4th ed.) New York: Wiley.
- Berger, J. (1991) Robust Bayesian analysis: Sensitivity to the prior. *J. Statist. Plann. Infer.*, **25**, 303–328.
- Berger, J. and Berliner, L. M. (1986) Robust Bayes and empirical Bayes analysis with  $\varepsilon$ -contaminated priors. *Ann. Statist.*, **14**, 461–486.
- Berger, J. O. and Wolpert, R. L. (1988) *The Likelihood Principle*. In *IMS Lecture Notes*, Ed. Gupta, S. S., Institute of Mathematical Statistics, vol. 6.

- Christensen, R., Johnson, W. O., Branscum, A. J. and Hanson, T. E. (2011) *Bayesian Ideas and Data Analysis*. Boca Raton: CRC Press.
- Cheney, W. (2001) *Analysis for Applied Mathematics*. New York: Springer.
- Evans, M. and Jang, G. H. (2011) Weak informativity and the information in one prior relative to another. *Statist. Sci.*, **26**, 423–439.
- Evans, M. and Moshonov, H. (2006) Checking for prior-data conflict. *Bayesian Anal.*, **1**, 893–914.
- Gelman, A., Jakulin, A., Pittau, M. G. and Su, Y. S. (2008) A weakly informative default prior distribution for logistic and other regression models. *Ann. Appl. Statist.*, **2**, 1360–1383.
- Giné, E. and Nickl, R. (2008) A simple adaptive estimator of the integrated square of a density. *Bernoulli*, **14**, 47–61.
- Grogan, W. and Wirth, W. (1981) A new American genus of predaceous midges related to palpomyia and bezzia (diptera: ceratopogonidae). *Proceedings of the Biological Society of Washington*, **94**, pp. 1279–1305.
- Hastie, T., Tibshirani, R. and Friedman, J. (2008) *Elements of Statistical Learning*. New York: Springer.
- Hoff, P. (2009) *A First Course in Bayesian Statistical Methods*. New York: Springer.
- Hunter, J. and Nachtergaele, B. (2005) *Applied Analysis*. London: World Scientific Publishing.
- Kyung, M., Gill, J., Ghosh, M. and Casella, G. (2010) Penalized regression, standard errors and Bayesian lassos. *Bayesian Anal.*, **5**, 369–412.
- Knight, K. (2000) *Mathematical Statistics*. Boca Raton: Chapman & Hall/CRC Press.
- Lavine, M. (1991) Sensitivity in Bayesian statistics: The prior and the likelihood. *J. Am. Statist. Assoc.*, **86** 396–399.
- Lenk, P. (2009) Simulation pseudo-bias correction to the harmonic mean estimator of integrated likelihoods. *J. Computnl Graph. Statist.*, **18**, 941–960.
- Lopes, H. F. and Tobias, J. L. (2011) Confronting prior convictions: On issues of prior sensitivity and likelihood robustness in Bayesian analysis. *Annual Rev. Econ.*, **3**, 107–131.
- Millman, R. S. and Parker, G. D. (1991) *Geometry: A Metric Approach with Models*. New York: Springer.
- Newton, M. A. and Raftery, A. E. (1994) Approximate Bayesian inference with the weighted likelihood Bootstrap (With Discussion). *J. R. Statist. Soc. B*, **56**, 3–26.
- Pajor, A. and Osiewalski, J. (2013) A note on Lenk’s correction of the harmonic mean estimator. *Central European Journal of Economic Modelling and Econometrics*, **5**, 271–275.
- Park, T. and Casella, G. (2008) The Bayesian lasso. *J. Am. Statist. Assoc.*, **103**, 681–686.
- Raftery, A. E., Newton, M. A., Satagopan, J. M. and Krivitsky, P. N. (2007) Estimating the integrated likelihood via posterior simulation using the harmonic mean identity. In *Bayesian Statistics*, Eds. Bernardo, J. M., Bayarri, M. J., Berger, J. O., Dawid, A. P., Heckerman, D., Smith, A. F. M. and West, M., Oxford University Press, vol. 8.
- Ramsey, J. O. and Silverman, B. W. (1997) *Functional Data Analysis*. New York: Springer-Verlag.
- Scheel, I., Green, P. J. and Rougier, J. C. (2011) A graphical diagnostic for identifying influential model choices in Bayesian hierarchical models. *Scandinavian J. Statist.*, **38**, 529–550.
- Shortle, J. F. and Mendel, M. B. (1996) The geometry of Bayesian inference, In *Bayesian Statistics*. eds. Bernardo, J. M., Berger, J. O., Dawid, A. P. and Smith, A. F. M., Oxford University Press, vol. 5, pp. 739–746.
- van der Vaart, A. W. (1998) *Asymptotic Statistics*. Cambridge: Cambridge University Press.
- Walter, G. and Augustin, T. (2009) Imprecision and prior-data conflict in generalized Bayesian inference. *J. Statist. Theo. Pract.*, **3**, 255–271.
- Wolpert, R. and Schmidler, S. (2012)  $\alpha$ -stable limit laws for harmonic mean estimators of marginal likelihoods. *Statist. Sinica*, **22**, 655–679.

Zhu, H., Ibrahim, J. G. and Tang, N. (2011) Bayesian influence analysis: A geometric approach. *Biometrika*, **98**, 307–323.